

Outrepasser les limites des techniques classiques de Prise d’Empreintes grâce aux Réseaux de Neurones

Javier Burroni¹ and Carlos Sarraute²

¹ Equipe de développement d’Impact, Core Security Technologies.

² Laboratoire de recherche de Core Security Technologies et Département de Mathématiques de l’Université de Buenos Aires.

Résumé Nous présentons la détection distante de systèmes d’exploitation comme un problème d’inférence : à partir d’une série d’observations (les réponses de la machine cible à un ensemble de tests), nous voulons inférer le type de système d’exploitation qui générerait ces observations avec une plus grande probabilité. Les techniques classiques utilisées pour réaliser cette analyse présentent plusieurs limitations. Pour outrepasser ces limites, nous proposons l’utilisation de Réseaux de Neurones et d’outils statistiques. Nous présenterons deux modules fonctionnels : un module qui utilise les points finaux DCE-RPC pour distinguer les versions de Windows, et un module qui utilise les signatures de Nmap pour distinguer les versions de systèmes Windows, Linux, Solaris, OpenBSD, FreeBSD et NetBSD. Nous expliquerons les détails de la topologie et du fonctionnement des réseaux de neurones utilisés, et du réglage fin de leurs paramètres. Finalement nous montrerons des résultats expérimentaux positifs.

1 Introduction

Le problème de la détection à distance du système d’exploitation, aussi appelé OS Fingerprinting (prise d’empreintes), est une étape cruciale d’un test de pénétration, puisque l’attaquant (professionnel de la sécurité, consultant ou hacker) a besoin de connaître l’OS de la machine cible afin de choisir les exploits qu’il va utiliser. La détection d’OS est accomplie en sniffant de façon passive des paquets réseau et en envoyant de façon active des paquets de test à la machine cible, pour étudier des variations spécifiques dans les réponses qui révèlent son système d’exploitation.

Les premières implémentations de prise d’empreintes étaient basées sur l’analyse des différences entre les implémentations de la pile TCP/IP. La génération suivante utilisa des données de la couche d’applications, tels que les endpoints DCE RPC. Bien que l’analyse porte sur plus d’information, quelque variation de l’algorithme consistant à chercher le point le plus proche (“best fit”) est toujours utilisée pour interpréter cette information. Cette stratégie a plusieurs défauts : elle ne va pas marcher dans des situations non standard, et ne permet

pas d'extraire les éléments clefs qui identifient de façon unique un système d'exploitation. Nous pensons que le prochain pas est de travailler sur les techniques utilisées pour analyser les données.

Notre nouvelle approche se base sur l'analyse de la composition de l'information relevée durant le processus d'identification du système pour découvrir les éléments clef et leurs relations. Pour implémenter cette approche, nous avons développé des outils qui utilisent des réseaux de neurones et des techniques des domaines de l'Intelligence Artificielle et les Statistiques. Ces outils ont été intégrés avec succès dans un software commercial (Core Impact).

2 DCE-RPC Endpoint mapper

2.1 Le service DCE-RPC nous informe

Dans les systèmes Windows, le service DCE (Distributed Computing Environment) RPC (Remote Procedure Call) reçoit les connexions envoyées au port 135 de la machine cible. En envoyant une requête RPC, il est possible de déterminer quels services ou programmes sont enregistrés dans la base de données du mapper d'endpoints RPC. La réponse inclut le UUID (Universal Unique Identifier) de chaque programme, le nom annoté, le protocole utilisé, l'adresse de réseau à laquelle le programme est lié, et le point final du programme (endpoint).

Il est possible de distinguer les versions, éditions et service packs de Windows selon la combinaison de point finaux fournie par le service DCE-RPC.

Par exemple, une machine Windows 2000 édition professionnelle service pack 0, le service RPC retourne 8 points finaux qui correspondent à 3 programmes :

```

uuid="5A7B91F8-FF00-11D0-A9B2-00C04FB6E6FC"
annotation="Messenger Service"
  protocol="ncalrpc"      endpoint="ntsvcs"      id="msgsvc.1"
  protocol="ncacn_np"     endpoint="\PIPE\ntsvcs" id="msgsvc.2"
  protocol="ncacn_np"     endpoint="\PIPE\scerpc" id="msgsvc.3"
  protocol="ncadg_ip_udp" endpoint=""        id="msgsvc.4"

uuid="1FF70682-0A51-30E8-076D-740BE8CEE98B"
  protocol="ncalrpc"      endpoint="LRPC"        id="mstask.1"
  protocol="ncacn_ip_tcp" endpoint=""            id="mstask.2"

uuid="378E52B0-C0A9-11CF-822D-00AA0051E40F"
  protocol="ncalrpc"      endpoint="LRPC"        id="mstask.3"
  protocol="ncacn_ip_tcp" endpoint=""            id="mstask.4"

```

2.2 Les réseaux de neurones entrent en jeu ...

Notre idée est de modeler la fonction qui fait correspondre les combinaisons de points finaux aux versions du système d'exploitation avec un réseau de neurones.

Plusieurs questions se posent :

- quel genre de réseau de neurones allons nous utiliser ?
- comment organiser les neurones ?
- comment faire correspondre les combinaisons de points finaux avec les neurones d'entrée du réseau ?
- comment entraîner le réseau ?

Le choix adopté est d'utiliser un réseau perceptron multicouches, plus précisément composé de 3 couches (nous indiquons entre parenthèses le nombre de neurones pour chaque couche qui résulta des tests de notre laboratoire, voir la figure 1 pour un schéma de la topologie du réseau).

1. couche d'entrée (avec 413 neurones) contient un neurone pour chaque UUID et un neurone pour chaque point final qui correspond à cet UUID. Suivant l'exemple précédent, nous avons un neurone pour le service Messenger et 4 neurones pour chaque endpoint associé à ce programme. Cela nous permet de répondre avec flexibilité à l'apparition d'un point final inconnu : nous retenons de toute façon l'information de l'UUID principal.
2. couche de neurones cachés (avec 42 neurones), où chaque neurone représente une combinaison des neurones d'entrée.
3. couche de sortie (avec 25 neurones), contient un neurone pour chaque version et édition de Windows (p.ex. Windows 2000 édition professionnelle), et un neurone pour chaque version et service pack de Windows (p.ex. Windows 2000 service pack 2). De cette manière le réseau peut distinguer l'édition et le service pack de façon indépendante : les erreurs dans une dimension n'affectent pas les erreurs dans l'autre dimension.

Qu'est-ce qu'un perceptron ? C'est l'unité de base du réseau, et en suivant l'analogie biologique, il joue le rôle d'un neurone qui reçoit de l'énergie à travers de ses connections synaptiques et la transmet à son tour selon une fonction d'activation.

Concrètement, chaque perceptron calcule sa valeur de sortie comme

$$v_{i,j} = f\left(\sum_{k=0}^n w_{i,j,k} \cdot x_k\right)$$

où $\{x_1 \dots x_n\}$ sont les entrées du neurone, $x_0 = -1$ est une entrée de biais fixe et f est une fonction d'activation non linéaire (nous utilisons tanh, la tangente hyperbolique). L'entraînement du réseau consiste à calculer les poids synaptiques $\{w_{i,j,0} \dots w_{i,j,n}\}$ pour chaque neurone (i, j) (i est la couche, j est la position du neurone dans la couche).

2.3 Entraînement par rétropropagation

L'entraînement se réalise par rétropropagation : pour la couche de sortie, à partir d'une sortie attendue $\{y_1 \dots y_m\}$ nous calculons (pour chaque neurone)

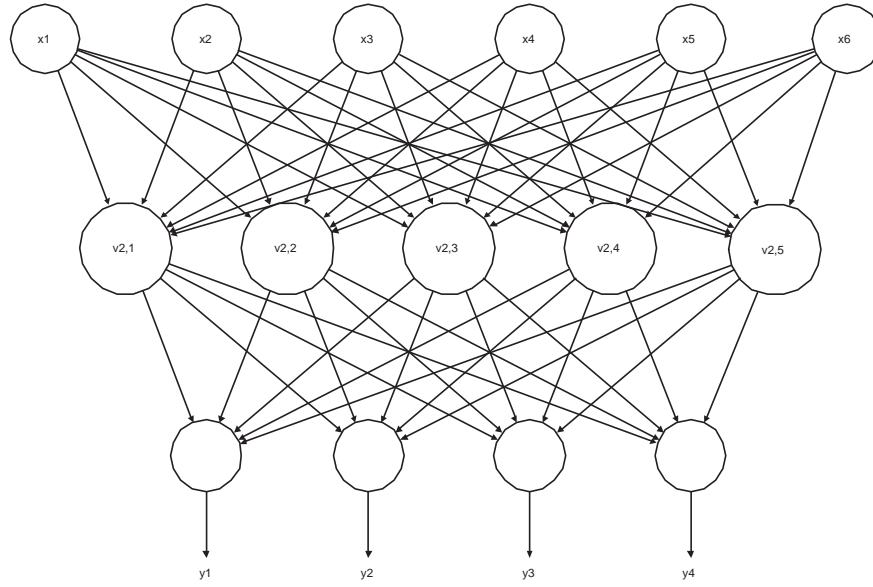


Figure 1. Réseau de Neurones Perceptron Multicouches

une estimation de l'erreur

$$\delta_{i,j} = f'(v_{i,j}) (y_j - v_{i,j})$$

Celle-ci est propagée aux couches précédentes par :

$$\delta_{i,j} = f'(v_{i,j}) \sum_k w_{i,j,k} \cdot \delta_{i+1,j}$$

où la somme est calculée sur tous les neurones connectés avec le neurone (i, j) .

Les nouveaux poids, au temps $t + 1$, sont

$$w_{t+1;i,j,k} = w_{t;i,j,k} + \Delta w_{t;i,j,k}$$

où Δw_t dépend d'un facteur de correction et aussi de la valeur de Δw_{t-1} multipliée par un momentum μ (cela donne aux modifications une sorte d'énergie cinétique) :

$$\Delta w_{t;i,j,k} = (\lambda \cdot \delta_{i+1,k} \cdot v_{i,j}) + \mu \cdot \Delta w_{t-1;i,j,k}$$

Le facteur de correction dépend des valeurs δ calculées et aussi d'un facteur d'apprentissage λ qui peut être ajusté pour accélérer la convergence du réseau.

Le type d'entraînement réalisé s'appelle apprentissage supervisé (basé sur des exemples). La rétropropagation part d'un jeu de données avec des entrées et des sorties recherchées. Une génération consiste à recalculer les poids synaptiques pour chaque paire d'entrée / sortie. L'entraînement complet requiert

10350 générations, ce qui prend quelques 14 heures d'exécution de code python. Etant donné que le design de la topologie du réseau est un processus assez artisanal (essai et erreur), qui requiert d'entraîner le réseau pour chaque variation de la topologie afin de voir si elle produit de meilleurs résultats, le temps d'exécution nous a particulièrement motivés à améliorer la vitesse de convergence de l'entraînement (problème que nous traitons dans la section 4).

Pour chaque génération du processus d'entraînement, les entrées sont réordonnées au hasard (pour que l'ordre des données n'affecte pas l'entraînement). En effet le réseau est susceptible d'apprendre n'importe quel aspect des entrées fournies, en particulier l'ordre dans lequel on lui présente les choses.

La table 1 montre une comparaison (de notre laboratoire) entre le vieux module DCE-RPC (qui utilise un algorithme de "best fit") et le nouveau module qui utilise un réseau de neurones pour analyser l'information.

Résultat	Ancien module DCE-RPC	DCE-RPC avec réseau de neurones
Concordance parfaite	6	7
Concordance partielle	8	14
Erreur	7	0
Pas de réponse	2	2

Table 1. Comparaison entre l'ancien module DCE-RPC et le nouveau module

La table 2 montre le résultat de l'exécution du module contre un Windows 2000 édition server sp1. Le système correct est reconnu avec un haut niveau de confiance.

3 Détection d'OS basée sur les signatures de Nmap

3.1 Richesse et faiblesse de Nmap

Nmap est un outil d'exploration réseau et un scanner de sécurité qui inclut une méthode de détection d'OS basée sur la réponse de la machine cible à une série de 9 tests. Nous décrivons brièvement dans le tableau 3 les paquets envoyés par chaque test, pour plus d'informations voir la page de Nmap [1].

Notre méthode utilise la base de signatures de Nmap. Une signature est un ensemble de règles décrivant comment une version / édition spécifique d'un système d'exploitation répond aux tests. Par exemple :

```
# Linux 2.6.0-test5 x86
Fingerprint Linux 2.6.0-test5 x86
Class Linux | Linux | 2.6.X | general purpose
TSeq(Class=RI%gcd=<6%SI=<2D3CFA0&>73C6B%IPID=Z%TS=1000HZ)
T1(DF=Y%W=16A0%ACK=S++%Flags=AS%Ops=MNNTNW)
```

```

Neural Network Output (close to 1 is better):
Windows NT4: 4.87480503763e-005
Editions:
    Enterprise Server: 0.00972694324639
    Server: -0.00963500026763
Service Packs:
    6: 0.00559659167371
    6a: -0.00846224120952
Windows 2000: 0.996048928128
Editions:
    Server: 0.977780526016
    Professional: 0.00868998746624
    Advanced Server: -0.00564873813703
Service Packs:
    4: -0.00505441088081
    2: -0.00285674134367
    3: -0.0093665583402
    0: -0.00320117552666
    1: 0.921351036343
Windows 2003: 0.00302898647853
Editions:
    Web Edition: 0.00128127138728
    Enterprise Edition: 0.00771786077082
    Standard Edition: -0.0077145024893
Service Packs:
    0: 0.000853988551952
Windows XP: 0.00605168045887
Editions:
    Professional: 0.00115635710749
    Home: 0.000408057333416
Service Packs:
    2: -0.00160404945542
    0: 0.00216065240615
    1: 0.000759109188052
Setting OS to Windows 2000 Server sp1
Setting architecture: i386

```

Table 2. Résultat du module DCE-RPC endpoint mapper

Test	envoyer paquet	à un port	avec les flags
T1	TCP	TCP ouvert	SYN, ECN-Echo
T2	TCP	TCP ouvert	no flags
T3	TCP	TCP ouvert	URG, PSH, SYN, FIN
T4	TCP	TCP ouvert	ACK
T5	TCP	TCP fermé	SYN
T6	TCP	TCP fermé	ACK
T7	TCP	TCP fermé	URG, PSH, FIN
PU	UDP	UDP fermé	
TSeq	TCP * 6	TCP ouvert	SYN

Table 3. Description des paquets envoyés

T2 (Resp=Y%DF=Y%W=0%ACK=S%Flags=AR%Ops=)
 T3 (Resp=Y%DF=Y%W=16A0%ACK=S++%Flags=AS%Ops=MNNTNW)
 T4 (DF=Y%W=0%ACK=0%Flags=R%Ops=)
 T5 (DF=Y%W=0%ACK=S++%Flags=AR%Ops=)
 T6 (DF=Y%W=0%ACK=0%Flags=R%Ops=)
 T7 (DF=Y%W=0%ACK=S++%Flags=AR%Ops=)
 PU (DF=N%TOS=C0%IPLen=164%RIPTL=148%RID=E%RIPCK=E%UCK=E%ULEN=134%DAT=E)

La base de Nmap contient 1684 signatures, ce qui veut dire que quelques 1684 versions / éditions de systèmes d'exploitation peuvent théoriquement être distinguées par cette méthode.

Nmap fonctionne en comparant la réponse d'une machine avec chaque signature de la base de données. Un score est assigné à chaque signature, calculé simplement comme le nombre de règles qui concordent divisé par le nombre de règles considérées (en effet, les signatures peuvent avoir différents nombres de règles, ou quelques champs peuvent manquer dans la réponse, dans ce cas les règles correspondantes ne sont pas prises en compte). C'est-à-dire que Nmap effectue une espèce de "best fit" basé sur une distance de Hamming, où tous les champs de la réponse ont le même poids.

Un des problèmes que présente cette méthode est le suivant : les systèmes d'exploitation rares (improbables) qui génèrent moins de réponses aux tests obtiennent un meilleur score ! (les règles qui concordent acquièrent un plus grand poids relatif). Par exemple, il arrive que Nmap détecte un Windows 2000 comme un Atari 2600 ou un HP-UX ... La richesse de la base de données devient alors une faiblesse !

3.2 Structure de Réseaux Hiérarchique

Si nous représentons symboliquement l'espace des systèmes d'exploitation comme un espace à 568 dimensions (nous verrons par la suite le pourquoi de ce nombre), les réponses possibles des différentes versions des systèmes inclus dans la base de données forme un nuage de points. Ce grand nuage est structuré de

façon particulière, puisque les familles de systèmes d'exploitation forment des clusters plus ou moins reconnaissables. La méthode de Nmap consiste, à partir de la réponse d'une machine, à chercher le point le plus proche (selon la distance de Hamming déjà mentionnée).

Notre approche consiste en premier lieu à filtrer les systèmes d'exploitation qui ne nous intéressent pas (toujours selon le point de vue de l'attaquant, par exemple les systèmes pour lesquels il n'a pas d'exploits). Dans notre implémentation, nous sommes intéressés par les familles Windows, Linux, Solaris, OpenBSD, NetBSD et FreeBSD. Ensuite nous utilisons la structure des familles de systèmes d'exploitation pour assigner la machine à une des 6 familles considérées.

Le résultat est un module qui utilise plusieurs réseaux de neurones organisés de façon hiérarchique :

1. premier pas, un réseau de neurones pour décider si l'OS est intéressant ou non.
2. deuxième pas, un réseau de neurones pour décider la famille de l'OS : Windows, Linux, Solaris, OpenBSD, FreeBSD, NetBSD.
3. dans le cas de Windows, nous utilisons le module DCE-RPC endpoint mapper pour raffiner la détection.
4. dans le cas de Linux, nous réalisons une analyse conditionnée (avec un autre réseau de neurones) pour décider la version du kernel.
5. dans le cas de Solaris et des BSD, nous réalisons une analyse conditionnée pour décider la version.

Nous utilisons un réseau de neurones différent pour chaque analyse. Nous avons ainsi 5 réseaux de neurones, et chacun requiert une topologie et un entraînement spécial.

3.3 Entrées du réseau de neurones

La première question à résoudre est : comment traduire la réponse d'une machine en entrées pour le réseau de neurones ? Nous assignons un ensemble de neurones d'entrée à chaque test.

Voici les détails pour les tests T1 ... T7 :

- un neurone pour le flag ACK.
- un neurone pour chaque réponse : S, S++, O.
- un neurone pour le flag DF.
- un neurone pour la réponse : yes/no.
- un neurone pour le champ *Flags*.
- un neurone pour chaque flag : ECE, URG, ACK, PSH, RST, SYN, FIN (en total 8 neurones).
- 10 groupes de 6 neurones pour le champ *Options*. Nous activons un seul neurone dans chaque groupe suivant l'option - EOL, MAXSEG, NOP, TIMESTAMP, WINDOW, ECHOED - en respectant l'ordre d'apparition (soit au total 60 neurones pour les options).

· un neurone pour le champ W (window size), qui a pour entrée une valeur hexadécimale.

Pour les flags ou les options, l'entrée est 1 ou -1 (présent ou absent). D'autres neurones ont une entrée numérique, comme le champ W (window size), le champ GCD (plus grand commun diviseur des numéros de séquence initiaux) ou les champs SI et VAL des réponses au test Tseq. Dans l'exemple d'un Linux 2.6.0, la réponse :

T3(Resp=Y%DF=Y%W=16A0%ACK=S++%Flags=AS%Ops=MNNTNW)

ACK	S	S++	O	DF	Yes	Flags	E	U	A	P	R	S	F	...
1	-1	1	-1	1	1	1	-1	-1	1	-1	-1	1	-1	...

Table 4. Résultat de la transformation (Linux 2.6.0)

se transforme en comme indiqué dans le tableau 4. De cette façon nous obtenons une couche d'entrée avec 568 dimensions, avec une certaine redondance. La redondance nous permet de traiter de façon flexible les réponses inconnues mais introduit aussi des problèmes de performance ! Nous verrons dans la section suivante comment résoudre ce problème (en réduisant le nombre de dimensions). Comme pour le module DCE-RPC, les réseaux de neurones sont composés de 3 couches. Par exemple le premier réseaux de neurones (le filtre de pertinence) contient : la couche d'entrée 96 neurones, la couche cachée 20 neurones, la couche de sortie 1 neurone.

3.4 Génération du jeu de données

Pour entraîner le réseau de neurones nous avons besoin d'entrées (réponses de machines) avec les sorties correspondantes (OS de la machine). Comme la base de signatures contient 1684 règles, nous estimons qu'une population de 15000 machines est nécessaire pour entraîner le réseau. Nous n'avons pas accès à une telle population ... et scanner l'Internet n'est pas une option !

La solution que nous avons adoptée est de générer les entrées par une simulation Monte Carlo. Pour chaque règle, nous générons des entrées correspondant à cette règle. Le nombre d'entrées dépend de la distribution empirique des OS basée sur des données statistiques. Quand la règle spécifie une constante, nous utilisons cette valeur, et quand la règle spécifie des options ou un intervalle de valeurs, nous choisissons une valeur en suivant une distribution uniforme.

4 Réduction des dimensions et entraînement

4.1 Matrice de corrélation

Lors du design de la topologie des réseaux, nous avons été généreux avec les entrées : 568 dimensions, avec une redondance importante. Une conséquence

est que la convergence de l'entraînement est lente, d'autant plus que le jeu de données est très grand. La solution à ce problème fut de réduire le nombre de dimensions. Cette analyse nous permet aussi de mieux comprendre les éléments importants des tests utilisés.

Le premier pas est de considérer chaque dimension d'entrée comme une variable aléatoire X_i ($1 \leq i \leq 568$). Les dimensions d'entrée ont des ordres de grandeur différents : les flags prennent comme valeur 1/-1 alors que le champ ISN (numéro de séquence initial) est un entier de 32 bits. Nous évitons au réseau de neurones d'avoir à apprendre à additionner correctement ces variables hétérogènes en normalisant les variables aléatoires (en soustrayant la moyenne μ et en divisant par l'écart type σ) :

$$\frac{X_i - \mu_i}{\sigma_i}$$

Puis nous calculons la matrice de corrélation R , dont les éléments sont :

$$R_{i,j} = \frac{E[(X_i - \mu_i)(X_j - \mu_j)]}{\sigma_i \sigma_j}$$

Le symbole E désigne l'espérance mathématique. Puisqu'après avoir normalisé les variables, $\mu_i = 0$ et $\sigma_i = 1$ pour tout i , la matrice de corrélation est simplement $R_{i,j} = E[X_i X_j]$.

La corrélation est une mesure de la dépendance statistique entre deux variables (une valeur proche de 1 ou -1 indique une plus forte dépendance). La dépendance linéaire entre des colonnes de R indique des variables dépendantes, dans ce cas nous en gardons une et éliminons les autres, puisqu'elles n'apportent pas d'information additionnelle. Les constantes ont une variance nulle et sont aussi éliminées par cette analyse.

Voyons le résultat dans le cas des systèmes OpenBSD. Nous reproduisons ci-dessous les extraits des signatures de deux OpenBSD différents, où les champs qui survivent à la réduction de la matrice de corrélation sont marqués en italiques.

```
Fingerprint OpenBSD 3.6 (i386)
Class OpenBSD | OpenBSD | 3.X | general purpose
T1(DF=N % W=4000 % ACK=S++ % Flags=AS % Ops=MNWNNT)
T2(Resp=N)
T3(Resp=N)
T4(DF=N % W=0 % ACK=0 %Flags=R % Ops=)
T5(DF=N % W=0 % ACK=S++ % Flags=AR % Ops=)
Fingerprint OpenBSD 2.2 - 2.3
Class OpenBSD | OpenBSD | 2.X | general purpose
T1(DF=N % W=402E % ACK=S++ % Flags=AS % Ops=MNWNNT)
T2(Resp=N)
T3(Resp=Y % DF=N % W=402E % ACK=S++ % Flags=AS % Ops=MNWNNT)
T4(DF=N % W=4000 % ACK=0 % Flags=R % Ops=)
T5(DF=N % W=0 % ACK=S++ % Flags=AR % Ops=)
```

Par exemple, pour le test T1, les seuls champs qui varient sont W et les deux premières options, les autres sont constants dans toutes les versions de

OpenBSD. Autre exemple, pour le test T4 seul W est susceptible de varier, et le test T5 n'apporte directement aucune information sur la version de OpenBSD examinée.

Index	Ancien index	Nom du champ
0	20	T1 : TCP OPT 1 EOL
1	26	T1 : TCP OPT 2 EOL
2	29	T1 : TCP OPT 2 TIMESTAMP
3	74	T1 : W FIELD
4	75	T2 : ACK FIELD
5	149	T2 : W FIELD
6	150	T3 : ACK FIELD
7	170	T3 : TCP OPT 1 EOL
8	179	T3 : TCP OPT 2 TIMESTAMP
9	224	T3 : W FIELD
10	227	T4 : SEQ S
11	299	T4 : W FIELD
12	377	T6 : SEQ S
13	452	T7 : SEQ S
14	525	TSeq : CLASS FIELD
15	526	TSeq : SEQ TD
16	528	TSeq : SEQ RI
17	529	TSeq : SEQ TR
18	532	TSeq : GCD FIELD
19	533	TSeq : IPID FIELD
20	535	TSeq : IPID SEQ BROKEN INCR
21	536	TSeq : IPID SEQ RPI
22	537	TSeq : IPID SEQ RD
23	540	TSeq : SI FIELD
24	543	TSeq : TS SEQ 2HZ
25	546	TSeq : TS SEQ UNSUPPORTED
26	555	PU : UCK RID RIPCK EQ
27	558	PU : UCK RID RIPCK ZERO
28	559	PU : UCK RID RIPCK EQ
29	560	PU : UCK RID RIPCK FAIL
30	563	PU : UCK RID RIPCK EQ
31	564	PU : UCK RID RIPCK FAIL
32	565	PU : RIPTL FIELD
33	566	PU : TOS FIELD

Table 5. Champs qui permettent de distinguer les OpenBSD

La table 5 montre la liste complète des champs qui servent à distinguer les différentes versions d'OpenBSD. Comme nous l'avons dit, le test T5 n'apparaît pas, alors que les tests Tseq et PU conservent de nombreuses variables, ce qui

nous montre que ces deux tests sont les plus discriminatifs au sein de la population OpenBSD.

4.2 Analyse en Composantes Principales

Une réduction ultérieure des données utilise l'Analyse en Composantes Principales (ACP). L'idée est de calculer une nouvelle base (ou système de coordonnées) de l'espace d'entrée, de telle manière que la majeure variance de toute projection du jeu de données dans un sous-espace de k dimensions, provient de projeter sur les k premiers vecteurs de cette base.

L'algorithme ACP consiste à :

- calculer les vecteurs propres et valeurs propres de R .
- trier les vecteurs par valeur propre décroissante.
- garder les k premiers vecteurs pour projeter les données.
- le paramètre k est choisi pour maintenir au moins 98% de la variance totale.

Après avoir réalisé l'ACP nous avons obtenu les topologies indiquée dans le tableau 6 pour les réseaux de neurones (la taille de la couche d'entrée originale était de 568 dans tous les cas).

Analyse	Couche d'entrée (après réduction de la matrice R)	Couche d'entrée (après ACP)	Couche cachée	Couche de sortie
Pertinence	204	96	20	1
Famille d'OS	145	66	20	6
Linux	100	41	18	8
Solaris	55	26	7	5
OpenBSD	34	23	4	3

Table 6. Topologies des réseaux de neurones

Pour conclure l'exemple de OpenBSD, à partir des 34 variables qui ont survécu à la réduction de la matrice de corrélation, il est possible de construire une nouvelle base de 23 vecteurs. Les coordonnées dans cette base sont les entrées du réseau, la couche cachée ne contient que 4 neurones et la couche de sortie 3 neurones (car nous distinguons 3 groupes de versions d'OpenBSD). Une fois que l'on sait qu'une machine est un OpenBSD, le problème de reconnaître la version est beaucoup plus simple et borné, et peut être accompli par un réseau de neurones de petite taille (plus efficient et rapide).

4.3 Taux d'apprentissage adaptatif

C'est une stratégie pour accélérer la convergence de l'entraînement. Le taux d'apprentissage est le paramètre λ qui intervient dans les formules d'apprentissage par rétropropagation.

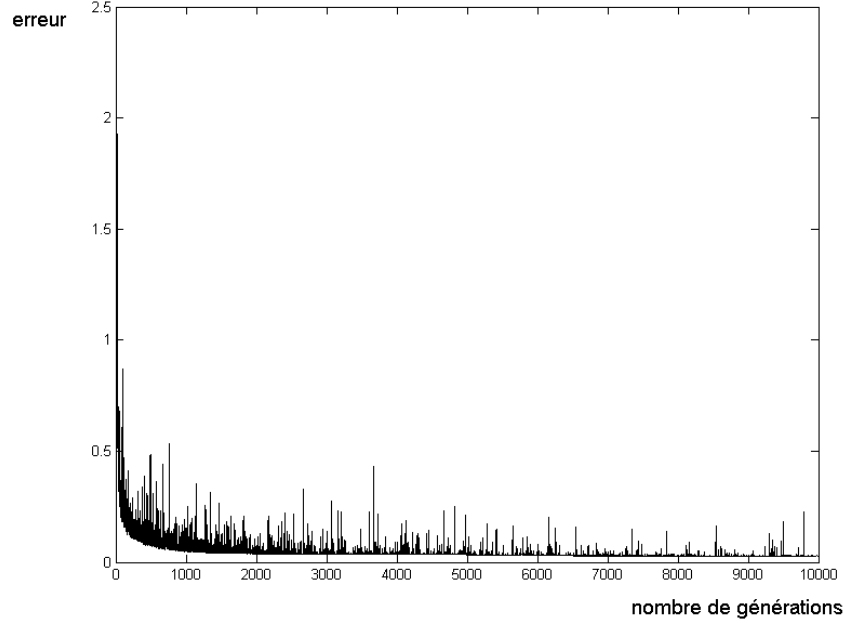


Figure 2. Taux d'apprentissage fixe

Etant donnée une sortie du réseau, nous pouvons calculer une estimation de l'erreur quadratique

$$\frac{\sum_{i=1}^n (y_i - v_i)^2}{n}$$

où y_i sont les sorties recherchées et v_i sont les sorties du réseau.

Après chaque génération (c'est-à-dire après avoir fait les calculs pour toutes les paires d'entrée / sortie), si l'erreur est plus grande, nous diminuons le taux d'apprentissage. Au contraire, si l'erreur est plus petite, alors nous augmentons le taux d'apprentissage. L'idée est de se déplacer plus rapidement si nous allons dans la direction correcte.

Voici deux graphiques (Figures 2 et 3) qui montrent l'évolution de l'erreur quadratique moyenne en fonction du nombre de générations pour chaque stratégie. Lorsque le taux d'apprentissage est fixe, l'erreur diminue et atteint des valeurs satisfaisantes après 4000 ou 5000 générations. L'erreur a une claire tendance à la baisse (les résultats sont bons) mais avec des pics irréguliers, dus à la nature probabiliste de l'entraînement du réseau.

En utilisant un taux d'apprentissage adaptatif, nous obtenons au début un comportement plus chaotique, avec des niveaux d'erreur plus hauts. Mais une fois que le système trouve la direction correcte, l'erreur chute rapidement pour

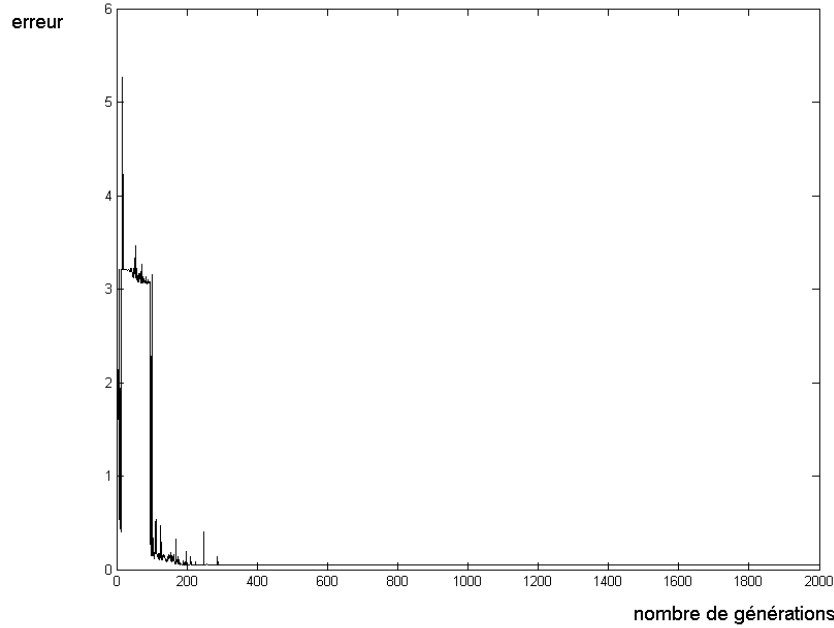


Figure 3. Taux d'apprentissage adaptatif

atteindre une valeur très faible et constante après 400 générations. Ces résultats sont clairement meilleurs et permettent d'accélérer l'entraînement des réseaux.

4.4 Entraînement par sous-ensembles du jeu de données

C'est une autre stratégie pour accélérer la convergence de l'entraînement. Elle consiste à entraîner le réseau avec plusieurs sous-ensembles plus petits du jeu de données. Ceci permet aussi de résoudre des problèmes de limitation d'espace, par exemple lorsque le jeu de données est trop grand pour être chargé entier en mémoire.

Pour estimer l'erreur commise après avoir utilisé un sous-ensemble, nous calculons une mesure d'adéquation G . Si la sortie est 0/1 :

$$G = 1 - (\text{Pr[faux positif]} + \text{Pr[faux négatif]})$$

Pour d'autres type de sorties, G est simplement :

$$G = 1 - \text{nombre d'erreurs/nombre de sorties.}$$

A nouveau, nous utilisons une stratégie de taux d'apprentissage adaptatif. Si la mesure d'adéquation G augmente, alors nous augmentons le taux d'apprentissage initial (pour chaque sous-ensemble).

Nous reproduisons ci-dessous le résultat de l'exécution du module contre une machine Solaris 8. Le système correct est reconnu avec précision.

```
Relevant / not relevant analysis
0.9999999999999789      relevant
```

```
Operating System analysis
-0.9999999999999434    Linux
0.99999999921394744    Solaris
-0.9999999999998057    OpenBSD
-0.99999964651426454    FreeBSD
-1.0000000000000000    NetBSD
-1.0000000000000000    Windows
```

```
Solaris version analysis
0.98172780325074482    Solaris 8
-0.99281382458335776    Solaris 9
-0.99357586906143880    Solaris 7
-0.99988378968003799    Solaris 2.X
-0.99999999977837983    Solaris 2.5.X
```

5 Conclusion et idées pour le futur

Dans ce travail, nous avons vu que l'une des principales limitations des techniques classiques de détection des systèmes d'exploitation réside dans l'analyse des données recueillies par les tests, basée sur quelque variation de l'algorithme de "best fit" (chercher le point le plus proche en fonction d'une distance de Hamming).

Nous avons vu comment générer et réunir l'information à analyser, comment homogénéiser les données (normaliser les variables d'entrée) et surtout comment dégager la structure des données d'entrée. C'est l'idée principale de notre approche, qui motive la décision d'utiliser des réseaux de neurones, de diviser l'analyse en plusieurs étapes hiérarchiques, et de réduire le nombre de dimensions d'entrée. Les résultats expérimentaux (de notre laboratoire) montrent que cette approche permet d'obtenir une reconnaissance plus fiable des systèmes d'exploitation.

De plus, la réduction de la matrice de corrélation et l'analyse en composantes principales, introduits en principe pour réduire le nombre de dimensions et améliorer la convergence de l'entraînement, nous donnent une méthode systématique pour analyser les réponses des machines aux stimuli envoyés. Cela nous a permis de dégager les éléments clefs des tests de Nmap, voir par exemple la table des champs permettant de distinguer les différentes versions d'OpenBSD. Une application de cette analyse serait d'optimiser les tests de Nmap pour générer moins de trafic. Une autre application plus ambitieuse serait de créer une base de données avec les réponses d'une population représentative de machines à une vaste batterie de tests (combinaisons de différents types de paquets, ports

et flags). Les mêmes méthodes d'analyse permettraient de dégager de cette vaste base de données les tests les plus discriminatifs pour la reconnaissance d'OS.

L'analyse que nous proposons peut aussi s'appliquer à d'autres méthodes de détection :

1. Xprobe2, d'Ofir Arkin, Fyodor & Meder Kydyraliev, qui base la détection sur des tests ICMP, SMB, SNMP.
2. Passive OS Identification (p0f) de Michal Zalewski, méthode qui a l'avantage de ne pas générer de trafic additionnel. C'est un défi intéressant, car l'analyse porte sur un volume de données plus important (tout le trafic sniffé), et requiert sans doute des méthodes plus dynamiques et évolutives.
3. Détection d'OS basée sur l'information fournie par le portmapper SUN RPC, permettant de distinguer des systèmes Sun, Linux et autres versions de System V.
4. Réunion d'information pour le versant client-side des tests de pénétration, en particulier pour détecter les versions d'applications. Par exemple détecter les Mail User Agents (MUA) tels que Outlook ou Thunderbird, en utilisant les Mail Headers.

Une autre idée pour le futur est d'ajouter du bruit et le filtre d'un firewall aux données étudiées. Ceci permettrait de détecter la présence d'un firewall, d'identifier différents firewalls et de faire des tests plus robustes.

Références

- [1] Fyodor, Remote OS detection via TCP/IP Stack FingerPrinting
<http://www.insecure.org/nmap/nmap-fingerprinting-article.html>
- [2] Principal Component Analysis
http://en.wikipedia.org/wiki/Principal_component_analysis
- [3] Projets de Corelabs
<http://www.coresecurity.com/corelabs/>
- [4] Christopher M. Bishop, *Neural Networks for Pattern Recognition*, Oxford University Press, 1995.
- [5] Timothy Masters, *Practical Neural Network Recipes in C++*, Academic Press, 1994.
- [6] Simon Haykin, *Neural Networks : A Comprehensive Foundation*, Prentice Hall, 2nd edition (1998).
- [7] Robert Hecht-Nielsen, *NeuroComputing*, Addison-Wesley, 1990.
- [8] Alberto. H. Landro, *Acerca de la probabilidad*, Ed. Coop - 2da Edición (2002)
- [9] W. Richard Stevens, *TCP/IP Illustrated*, Addison-Wesley Professional, 1993.